

PATENT ABSTRACTS OF JAPAN

(11) Publication number : 07-271829
(43) Date of publication of application : 20. 10. 1995

(51) Int. Cl. G06F 17/30
G06K 9/00

(21) Application number : 07-064381 (71) Applicant : XEROX CORP
(22) Date of filing : 23. 03. 1995 (72) Inventor : SPITZ A LAWRENCE
DIAS ANTONIO P

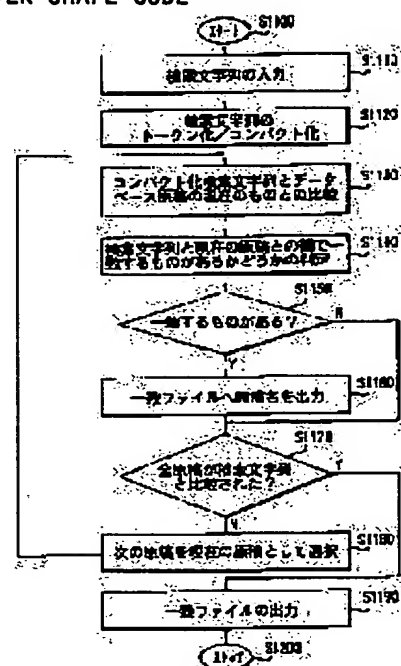
(30) Priority
Priority number : 94 220926 Priority date : 31. 03. 1994 Priority country : US

(54) MATCHING METHOD FOR TEXT IMAGE AND ORIGINAL USING CHARACTER SHAPE CODE

(57) Abstract:

PURPOSE: To accurately match a text image and a text original with each other by a compacted character shape code.

CONSTITUTION: The first method for performing accurate or non-strict matching of an original stored in an original data base includes a processing for converting a data base original into a compacted form made into a token. A retrieval character string or a retrieval original is converted into the compacted form made into the token and compared. Whether or not a test character string is present inside the data base original or whether or not the data base original corresponds to a test original is judged. A second, method for performing the non-strict matching of the test original and the data base original includes the processing for generating one or plural floating point number sets of the respective data base original and test original. The floating point number set of a data base is compared with the floating point number set of the test original and a matching degree is decided. A threshold value is used.



LEGAL STATUS

[Date of request for examination]
[Date of sending the examiner's decision of rejection]
[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]
[Date of final disposal for application]
[Patent number]
[Date of registration]
[Number of appeal against examiner's decision of rejection]
[Date of requesting appeal against examiner's decision of rejection]
[Date of extinction of right]

Copyright (C) ; 1998, 2003 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平7-271829

(43) 公開日 平成7年(1995)10月20日

(51) Int.Cl. ⁶	識別記号	庁内整理番号	F I	技術表示箇所
G 0 6 F 17/30				
G 0 6 K 9/00	S	9289-5L		
		9194-5L	G 0 6 F 15/ 403	3 5 0 C
		9194-5L	15/ 40	3 7 0 B

審査請求 未請求 請求項の数 1 O L (全 11 頁)

(21) 出願番号 特願平7-64381

(22) 出願日 平成7年(1995)3月23日

(31) 優先権主張番号 08/220926

(32) 優先日 1994年3月31日

(33) 優先権主張国 米国 (US)

(71) 出願人 590000798

ゼロックス コーポレーション

XEROX CORPORATION

アメリカ合衆国 ニューヨーク州 14644

ロチェスター ゼロックス スクエア

(番地なし)

(72) 発明者 エイ ローレンス スピッツ

アメリカ合衆国 カリフォルニア州

94303 パロ アルト サウザンブトン

ドライヴ 821

(74) 代理人 弁理士 中村 稔 (外6名)

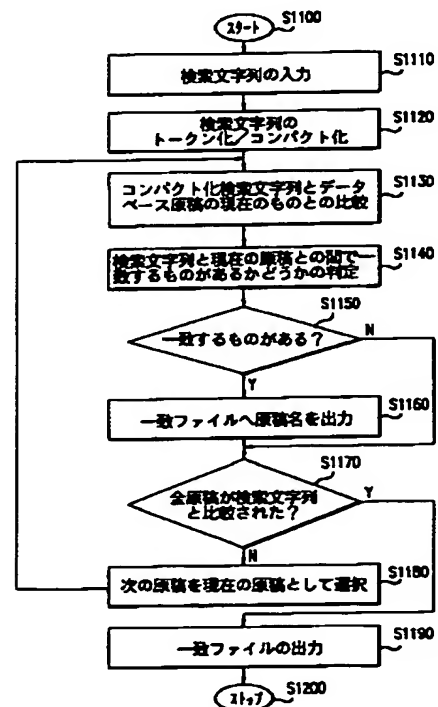
最終頁に続く

(54) 【発明の名称】 文字形状コードを用いたテキスト画像と原稿とのマッチング方法

(57) 【要約】

【目的】 コンパクト化された文字形状コードで、テキスト画像とテキスト原稿との正確なマッチングを行う。

【構成】 原稿データベースに記憶された原稿の正確または厳密でないマッチングを行う第1の方法は、データベース原稿をコンパクト化されトークン化された形態に変換する処理を含む。検索文字列または検索原稿はコンパクト化トークン化形態に変換され、比較される。テスト文字列がデータベース原稿内にあるかどうか、またはデータベース原稿がテスト原稿と対応するかどうか判断される。テスト原稿とデータベース原稿との厳密でないマッチングを行う第2の方法は、各データベース原稿およびテスト原稿の一または複数の浮動小数点数セットを生成する処理を含む。データベースの浮動小数点数セットはテスト原稿の浮動小数点数セットと比較され、一致度が決定される。しきい値が用いられる。



【特許請求の範囲】

【請求項 1】 テスト・ストリングを複数の原稿とマッチングするための方法であって、
前記複数の原稿のそれぞれを文字形状コード表記に変換し、
前記複数の原稿のそれぞれに対して、前記文字形状コード表記を、縮小されたバイナリ・データ・ストリングに変換し、
前記複数の原稿のそれぞれの縮小されたバイナリ・データ・ストリングをデータベースに格納し、
テストストリングを入力し、
前記テスト・ストリングを文字形状コード表記に変換し、
前記テスト・ストリングの前記文字形状コード表記を縮小されたバイナリ・データ・ストリングに変換し、
前記テスト・ストリングの前記縮小されたバイナリ・データ・ストリングを前記複数の原稿のそれぞれの縮小されたバイナリ・データ・ストリングとマッチングし、
マッチングした原稿のリストを出力する、
方法。

【発明の詳細な説明】

【0001】

【技術分野】 本発明は、テキスト画像（テキスト・イメージ）を原稿（ドキュメント）と比較するための方法、およびテキスト画像と原稿の文字形状コード表記の比較に基づき、原稿データベース内に原稿を発見するための方法に関する。特に、本発明は、テキスト画像および原稿の文字の全体の物理的形状を表す、制限された個数の文字コードを使用して、テキスト画像と原稿を文字形状コード表記に変換することに関する。

【0002】

【発明の背景】 この技術分野における周知の従来システムにおいては、テキスト画像をテキスト原稿と比較する前に、テキスト画像は、まず、光学式文字読取り（OCR: Optical Character Recognition）技術を用いて、分離したテキスト原稿に変換されなければならない。しかしながら、光学式文字読取り過程においては、一般に、認識テキスト原稿内に置換、消滅および挿入のような多くのエラーを有するものが生成される。このように、光学的に認識されたテキスト原稿と電子テキスト原稿との正確な照合を行うことは一般に不可能である。したがって、認識されたテキスト原稿を正確なデータベース・テキスト原稿と正しく照合することは、不可能でない場合であっても、困難であることが多い。

【0003】

【発明の概要】 本発明は、文字形状でコード化されたワード・トークンを含むテキスト画像を、コンパクト化された文字形状コード表記に変換するための方法を提供する。

【0004】 また、本発明は、コンパクト化された文字

形状コードで、テキスト画像とテキスト原稿との正確な照合を行う方法を提供する。

【0005】 さらに、本発明は、トークン化されたテキスト画像およびトークン化されたテキスト原稿のそれぞれから生成された複数の浮動小数点数を生成することにより、トークン化されたテキスト画像とテキスト原稿との厳密でない照合を行う方法を提供する。続いて、これらの浮動小数点数は比較され、テキスト画像とテキスト原稿とを厳密でない照合が行われる。

10 【0006】 本発明のこれらの目的および他の目的を達成するために、第 1 の好ましい実施例においては、テキスト画像は、まず走査され、続いてトークン化されたテキスト画像に変換される。この変換には、出願係属中の特許出願第 6-70296 号（参照のためここに引用）に示されている自動文字形状分類方法が用いられる。

【0007】 電子的に記憶された（すなわち、文字コード・テキストとして記憶された）テキスト原稿のライブラリは、この方法を用いて、トークン化された電子テキスト原稿（以下「トークン化電子テキスト原稿」ともいう）のライブラリに予め変換される。続いて、このトークン化電子テキスト原稿は、本発明のコンパクト化方法の第 1 の好ましい実施例を用いて、コンパクト化されたトークン化電子テキスト原稿に変換される。同様にして、テキスト画像がトークン化されたテキスト画像（以下「トークン化テキスト画像」ともいう）に変換された後に、これは、本発明のコンパクト化方法の第 1 の好ましい実施例を用いて、コンパクト化されたトークン化テキスト画像に変換される。続いて、このコンパクト化されたトークン化テキスト画像は、データベースに記憶された原稿のライブラリと比較され、正確に一致（マッチ）するものが決定される。

【0008】 第 2 の好ましい実施例においては、データベースは、コンパクト化されたトークン化電子テキスト原稿ではなく、コンパクト化されたトークン化テキスト画像のライブラリを含んでいる。電子テキスト文字列（電子テキスト・ストリング）は、コンパクト化され、トークン化された形態に変換され、テキスト画像のデータベース・ライブラリを検索するために使用され、そのテキスト文字列と正確に一致するものを含むテキスト画像が突き止められる。

【0009】 第 3 の好ましい実施例においては、テキスト画像および電子テキスト原稿は、一または二以上の浮動小数点数からなるセットによってそれぞれ表される。テスト原稿（テキスト画像または電子テキスト原稿のいずれか）の浮動小数点数のセットは、テキスト画像または電子テキスト原稿の複数の他のものを表す浮動小数点のセットのデータベースと比較される。データベース原稿の浮動小数点数のセットは、テスト原稿のセットと比較され、このテスト原稿と一致するテキスト画像または原稿の厳密でない識別が行われる。

【0010】

【実施例】テキスト画像（テキスト・イメージ）および電子テキスト原稿（電子テキスト・ドキュメント）が以下に示す方法を用いてコンパクト化される前に、このテキスト原稿およびテキスト画像は、まず、電子テキスト原稿の場合には文字コードからトークン化された文字形状コード表記に、テキスト画像の場合にはビットマップからトークン化された文字形状コード表記にそれぞれ変換されなければならない。

【0011】図1は、文字形状コード認識システムを示している。このシステムは、電荷結合素子（CCD）等を有するスキャナ110を備えている。スキャナ110は、図2に示すような原稿（ドキュメント）100を走査する。原稿100は画像（イメージ）102を有する。画像102は、ヨーロッパ・スクリプト・タイプのテキスト文字列を有する。スキャナ110は、オリジナルの原稿100の画像102を含む複数のピクセルの位置および画像濃度を表すデジタル・データ信号を出力する。

【0012】このデジタル・データ信号はメモリ112に送られる。メモリ112では、デジタル・データ信号が一時的または無期限に記憶される。デジタル・データ信号は、メモリ112から出力されると、汎用デジタル・コンピュータ114に入力される。コンピュータ114に入力されると、まず、デジタル・データ信号は、画像102の非テキスト部分を除去してテキスト部分104を残すことにより整理（クリーン・アップ）される。さらに、歪み等のデジタル・データ信号上のスキャナ生成物が補正される。整理されたデジタル・データ信号は、メモリ112に再び記憶されるか、または、コンピュータ114のメモリに記憶される。選択的に、スキャナ110がスキャナ生成物の除去のような前処理を提供することができる。

【0013】図1に示すように、本発明の汎用デジタル・コンピュータ114は、メモリ22および入出力回路24を備えている。メモリ22には、制御プログラムが記憶される。入出力回路24は、メモリ112からのデジタル・データ信号を入力し、画像102のテキスト部分104または一致（マッチ）する原稿の決定されたワード・トークンを表す信号を出力する。汎用コンピュータ114は、画像メモリ26、結合コンポーネント生成手段28、空間的特徴決定手段30、文字形状分類手段32、およびコンパクト手段34も備えている。画像メモリ26には、デジタル・データ信号が記憶される。結合コンポーネント生成手段28は、デジタル・データ信号から、結合したコンポーネント（結合コンポーネント）を生成する。空間的特徴決定手段30は、テキスト部104のライン（行）、ワードおよび文字セルの座標、ならびに各文字セル内の各結合コンポーネントの位置を決定する。文字形状分類手段32は、文字セルを、抽象化された文字形状コードに変換する。コンパクト手段34は、決定されたワード・トークン

に基づいて、テキスト部104を、コンパクト化され、かつ、トークン化された形に変換する。制御プログラムを記憶するメモリ22は、ROM22aまたはRAM22bのいずれを含んでいてもよい。

【0014】分類手段32の好ましい実施例においては、分類手段32は、結合コンポーネント計数手段320、トップ位置検出手段322、ボトム位置検出手段324、結合コンポーネント・サイジング手段326、ライン位置記憶手段328および比較手段330を備えている。結合コンポーネント計数手段320は、現在の文字セル内における結合コンポーネントの個数を決定する。トップ位置検出手段322は、現在の文字セル内の少なくとも一つの結合コンポーネントのトップの位置を突き止める。ボトム位置検出手段324は、現在の文字セルの少なくとも一つの結合コンポーネントのボトムの位置を突き止める。結合コンポーネント・サイジング手段326は、現在の文字セルの少なくとも一つの結合コンポーネントの高さと幅を決定する。ライン位置記憶手段328は、現在の文字セルを含んだラインの基底ラインおよびxラインの少なくとも一つのものを記憶する。比較手段330は、結合コンポーネントのトップ位置、結合コンポーネントのボトム位置および結合コンポーネントの高さの少なくとも一つを、基本ラインの位置、xラインの位置および結合コンポーネントの幅の少なくとも一つと比較する。

【0015】コンパクト化を行う方法の第1の好ましい実施例と関連した選択的な実施例においては、前記分類手段はギャップ決定手段332も含んでいる。ギャップ決定手段332は、結合コンポーネントの右側部分にギャップが存在するかどうかを判定する。もちろん、結合コンポーネント生成手段28、空間的特徴決定手段30、分類手段32および比較手段34の各機能および対応する手段を、独立した手段によって実現することもできるし、このような構造は、上述した本発明の好ましい実施例に相当することが分かる。コンパレータ36は、コンパクト化されたテキスト文字列または原稿をデータベースと比較する。また、コンパレータ36は、第4の好ましい実施例の浮動小数点の値のセット（集合）も比較する。

【0016】処理において、図2に示すような画像102を含む原稿100は、スキャナ110の上または内部に置かれ、走査される。そして、シリアルまたはパラレルのデジタル・データ信号が生成される。このデジタル・データ信号は複数の信号部分を含んでいる。各部分はオリジナル画像102に対応したピクセルを表す。画像102の各ピクセルは、画像102における位置および画像濃度を有する。したがって、デジタル・データ信号の各信号部分は、対応するピクセルの位置および画像濃度を表すデータを含んでいる。

【0017】続いて、スキャナ110から出力されたデジタル・データ信号は、メモリ112に記憶される。メモリ112は、RAM、フラッシュ・メモリ、ディスク・メ

メモリ等で構成することができる。メモリ112のタイプに関係なく、デジタル・データ信号は、各信号部分に含まれる位置および画像濃度データに応じてメモリ112に記憶される。もちろん、デジタル・データ信号を、この中継を行うメモリ112に記憶することなく、汎用デジタル・コンピュータ114に直接入力することもできる。選択的に、汎用デジタル・コンピュータ114にメモリ112を組み込むこともできる。メモリ112は、走査された画像102の長期記憶として使用される場合もある。

【0018】オペレータがスキャナ110への原稿の入力を完了するか、またはそうでなければ、システムが、画像102を表すデジタル・データ信号を文字形状コード・シンボルに変換すべきと判断すると、画像102を表すデジタル・データ信号はメモリ112から汎用コンピュータ114へ出力される。もちろん、特殊な用途のデジタル・コンピュータまたはハードウェア・ロジック回路を、汎用デジタル・コンピュータ114に代わって使用することもできる。

【0019】メモリ112に記憶されたデジタル画像データ信号は、汎用コンピュータ114に出力され、入出力手段24を介して画像メモリ26に入力される。デジタル・データ信号の画像メモリ26への記憶が完了すると、続いて、結合コンポーネント生成手段28がデジタル・データ信号を処理する。結合コンポーネント生成手段28は、画像102を表すデジタル・データ信号を複数の結合コンポーネントに分解する。各結合コンポーネントは、単一のラインの一または複数の信号部分を備えている。各結合コンポーネントは、オリジナル画像102におけるある最小の画像濃度を有し、かつ、連続した経路を形成するピクセルに対応した信号部分を含んでいる。一般に、各スクリプト文字は“Fuji”の“F”のように一つの結合コンポーネントに対応するものもあり、また、“Fuji”の“j”または“i”のように二以上の結合コンポーネントに対応するものもある。

【0020】結合コンポーネント生成手段28が、画像102の複数の結合コンポーネントをデジタル・データ信号から生成すると、画像102に対応したデジタル・データ信号、ならびに結合コンポーネント生成手段28によって生成された結合コンポーネントのリストおよびそれらの位置は、画像メモリ26に記憶され、空間的特徴決定手段30に出力される。

【0021】空間的特徴決定手段30は、ラインの位置、ワードのスペースおよび文字セルのようなテキスト部分の空間的特徴を決定する。各文字セルは、隣接するスペース間のライン内の垂直方向に並んだ結合コンポーネントを含んでいる。例えば、“Fuji”の文字“i”および“j”は、2つの独立した結合コンポーネントからそれぞれ形成されている。空間的特徴決定手段30は、一ラインの垂直方向に並んだ全ての結合コンポーネントを

一つの文字セルにグループ化する。結合コンポーネントの生成および結合コンポーネントからテキスト部分104の空間的特徴の決定を行う一つの好ましい方法および装置は、米国特許出願第07/047,514号に示されている。この米国出願は、本特許出願と同一の出願人により出願されたものであり、参照のためにここに引用される。

【0022】続いて、結合コンポーネントおよび文字セルのリストは、空間的特徴決定手段30によって出力され、文字形状コード分類手段32に与えられる。文字形状コード分類手段32は、文字セル内の結合コンポーネントまたは結合コンポーネント群を、文字セル内の結合コンポーネントの個数および位置に基づいて複数の抽象的な文字コード（抽象文字コード）の一つに変換する。図3は、好ましい文字コード・リストおよび各コードに対応する文字を示している。図3に示すように、13個の異なる抽象文字形状コードが使用される。しかしながら、第1および第2の好ましい実施例においては、これらの13個の文字形状コードはコード“A, U, i, x, g, j”に限定される。さらに、“スペース”および“CR”の2つの追加コードが使用される。コード“スペース”はワード間のスペースを示す。コード“CR”（キャリッジ・リターン）はラインの終了を示す。各抽象文字形状コードは、文字セル内の独立した結合コンポーネントの個数、各文字セルの独立した結合コンポーネント間の相対的な位置、および文字セル内の結合コンポーネントの位置に基づいて、一または複数の文字を代表するものである。

【0023】上述したトークン化を行うシステムの処理を簡単化したフローチャートが図4に示されている。ステップS100において、システムは処理を開始する。原稿がステップS110で走査され、デジタル・データ信号が生成される。続いて、ステップS120において、デジタル画像データ信号は整理（クリーンアップ）される。この整理は、所望の前処理アルゴリズムをこのデジタル画像データに適用することによって行われる。ステップS130において、デジタル画像データ信号の結合コンポーネントが識別され、ステップS140において、文字セルが決定される。ステップS150において、各文字セルの文字タイプ分類が決定される。ステップS160において、文字コードがともにグループ化され、ワード間およびワードの内部スペースに基づいてトークンが形成される。

【0024】図3に示すコード化を実行する決定木が図5に示されている。図5に示すように、1つの結合コンポーネントを有する文字セル用の7つの抽象文字コード、2つの結合コンポーネントを有する文字セル用の5つの抽象文字コード、および3つの結合コンポーネントを有する文字セル用の1つの抽象文字コードがある。

【0025】本発明の方法の好ましい実施例は図5に示す決定木を実行する。ステップS300において、分類手段32は、まず、現在の文字セル内の結合コンポーネントの

個数を決定する。本発明の好ましい実施例において、文字形状コード分類手段32は、テキスト部分104の各文字セルをセルごとに処理する。

【0026】本発明の方法および装置は統計的に強固な全体の特徴を分析するので、この方法および装置は非常に質の悪い印刷および（または）走査された原稿であっても処理でき、歪みが容易に生じる従来のOCR技術と対して、優れた特徴を分析できる。したがって、デジタル・データ信号またはこのデジタル・データ信号から生成された結合コンポーネントに対して、原稿の全ての文字を完全に表すことは必要ない。むしろ、本発明はよく起こる走査エラーに耐えることができる。すなわち、本発明は、例えば、1つの結合コンポーネント文字が2つまたはそれ以上の結合コンポーネントに分離したり、2つまたはそれ以上の分離した結合コンポーネントが1つの結合コンポーネントに合体したり、結合コンポーネントがライン上に間違っ

て置かれたりするようなエラーに対して強い。また、本発明の方法および装置は、歪みおよび（または）ねじれを有する画像を分析するときにも強さを発揮する。

【0027】ステップS300において、分類手段32が、文字セルが一つだけの結合コンポーネントを有すると判断すると、次に、分類手段32は、その結合コンポーネントのトップ位置が現在のラインのxライン（ミーンライン）位置より上にあり、かつ、ボトム位置がベースライン（並び線）より上にあるかどうかを判定する。ライン位置および結合コンポーネントの位置は、最上部の位置またはその近傍および最左端の位置またはその近傍にある基準位置から測定され、下方および右方がそれぞれ正になるように測定されることが分かる。

【0028】ステップS310が肯定的ならば、分類手段32は、ステップS320において文字セルをアポストロフィ（省略符号）に変換する。一方、ステップS310が否定的ならば、分類手段32はステップS330に進む。ステップS330において、分類手段32は、結合コンポーネントのトップ位置がxライン位置より上にあり、かつ、結合コンポーネントのボトムがベースライン位置上またはそれより下にあるかどうかを判定する。ステップS330が肯定的ならば、分類手段32はステップS340において文字セルを“A”に変換する。“A”は、図3に示すように、全ての大文字、全ての数字、アセンダを有する小文字、および一般に垂直方向（縦方向）を向いた句読点（パンクチュエーション・マーク）の全てを代表するものである。

【0029】ステップS330が否定的ならば、分類手段32はステップS350に進む。ステップS350において、分類手段32は、結合コンポーネントのトップがxライン位置より下にあり、かつ、結合コンポーネントのボトムがベースライン位置よりも上にあるかどうかを判定する。ステップS350が肯定的ならば、分類手段32はステップS360において文字セルをハイフンに変換する。

【0030】ステップS350が否定的ならば、分類手段32はステップS370に進む。ステップS370において、分類手段32は、結合コンポーネントのトップ位置がxライン位置より下にあり、かつ、結合コンポーネントのボトム位置がベースライン位置より下にあるかどうかを判定する。ステップS370が肯定的ならば、分類手段32は、ステップS380において文字セルをカンマに変換する。ステップS370が否定的ならば、分類手段32はステップS390に進む。ステップS390において、分類手段32は、結合コンポーネントのトップ位置がxライン位置の下にあるかどうかを判定する。ステップS390が肯定的ならば、分類手段32はステップS400において文字セルをピリオドに変換する。

【0031】ステップS390が否定的ならば、分類手段32はステップS410に進む。ステップS410において、分類手段32は、結合コンポーネントのボトム位置がベースライン位置より下にあるかどうかを判定する。ステップS410が肯定的ならば、分類手段32はステップS420でその文字セルを“g”に変換する。コード“g”は、図3に示すように、ディセンダを有する任意の小文字を代表するのである。

【0032】ステップS410が否定的ならば、分類手段32はステップS430に進む。ステップS430において、分類手段32は結合コンポーネントがアセンダまたはディセンダのいずれも有しない小文字であると仮定し、その結合コンポーネントを“x”に変換する。続いて、ステップS430、またはステップS320、S340、S360、S380、S400およびS420に続いて、分類手段32は次の文字セルを現在の文字セルとして選択して、ステップS300に戻る。

【0033】一方、ステップS300において、分類手段32が、現在の文字セル内に2つの結合コンポーネントがあると判断すると、分類手段32はステップS440に進む。ステップS440において、分類手段32は上部にある結合コンポーネントの高さが上部にある結合コンポーネントの幅の3倍よりも大きいかどうかを判定する。結合コンポーネントの高さとは、そのトップ位置とボトム位置との間の差であり、結合コンポーネントの幅はその右端位置と左端位置との間の差である。ステップS440が肯定的ならば、分類手段32はステップS450に進む。ステップS450において、分類手段32はその文字セルを感嘆符（!）に変換する。

【0034】ステップS440が否定的ならば、分類手段32はステップS460に進む。ステップS460において、分類手段32は上部にある結合コンポーネントのトップ位置がxライン位置より上にあり、かつ、下部にある結合コンポーネントのボトム位置がベースライン位置より下にあるかどうかを判定する。ステップS460が肯定的ならば、分類手段32はステップS470においてその文字セルを“j”に変換する。“j”は、xラインよりも上に延びる分離した結合コンポーネントおよびベースラインよりも下に

延びる分離した結合コンポーネントを有する任意の小文字を代表するものである。

【0035】ステップS460が否定的ならば、分類手段32はステップS480に進む。ステップS480において、分類手段32は、上部にある結合コンポーネントのトップ位置がxライン位置より上にあり、かつ、ボトム位置がベースライン位置より下にないかどうかを判定する。ステップS480が肯定的ならば、分類手段32はステップS490においてその文字セルを“i”に変換する。“i”は、図3に示すように、xライン位置より上に延びる分離した結合コンポーネントと、ベースライン位置より下には延びていない分離した結合コンポーネントとを有する任意の小文字を代表するものである。

【0036】ステップS480が否定的ならば、分類手段32はステップS500に進む。ステップS500において、分類手段32は、上部および下部にある結合コンポーネントの双方がそれらの高さの3倍の幅を有するかどうかを判定する。ステップS500が肯定的ならば、分類手段32はステップS510においてその文字セルを“=”に変換する。ステップS500が否定的ならば、分類手段32はその文字セルが“:”に変換されるべきと仮定して、ステップS520において、その文字セルはそうに変換される。ステップS520、ならびにステップS450、ステップS470、ステップS490およびステップS510の後、分類手段32は次の文字セルを現在の文字セルとして選択して、ステップS300に進む。

【0037】一方、分類手段32が、ステップS300において、現在の文字セル内に3つの結合コンポーネントがあると判断すると、分類手段32はステップS530に進む。ステップS530において、分類手段32は、文字セルがウムラウト符号を有する大文字または小文字を表すものと仮定して、これにより、図5に示すように、その文字セルを“U”に変換する。続いて、分類手段32は、次の文字セルを現在の文字セルとして選択し、ステップS300に進む。一方、次の文字セルがないならば、分類手段32はテキスト部分104の分類を終了し、文字セルの代わりに抽象文字コードのリストを画像メモリ26に出力する。このようにして、図2のテキスト画像は、図6に示す文字形状コード表記に変換される。

【0038】一または二以上のコード化された文字のワード・トークン・リストは、コンパクト化手段34に与えられる。コンパクト化手段34は、テキスト部分104を表すワード・トークンのリストを入力し、続いて“A、U、I、X、G”を除くトークンの全てをテキスト画像から消去して、コンパクト化され、かつ、トークン化されたテキスト部分104の形態を生成する。もちろん、ASCII (アスキー) のような文字コード体系ですでに表されている電子テキスト原稿の場合には、コントローラ114の文字コード変換器36は、一般に、ASCIIコードを文字形状コードに直接変換することができる。

【0039】トークン化されたテキスト画像（以下「トークン化テキスト画像」ともいう）およびトークン化された電子テキスト原稿（以下「トークン化電子テキスト原稿」ともいう）が生成されると、これらのトークン化テキスト画像およびトークン化電子テキスト原稿は、コンパクト化手段34によって、コンパクト化されたトークン化テキスト画像および原稿に変換される。この変換において、図7に示すように、コンパクト化方法の第1の好ましい実施例が用いられる。第1の好ましい実施例において、ステップS1000で、変換プロセスがスタートする。ステップS1010において、トークン化テキスト原稿は文字形状コードのフル・セットから文字形状コードの縮小セットに縮小される。これらの縮小された文字形状コード(A、U、g、i、j、x、[スペース]および[CR])は、実際の原稿に使用されるワード・トークンの大多数に使用されるコードを表す。可能な文字形状コードをこれらの8つの縮小された文字形状コードに限定することにより、文字形状コードを、8つの異なる3ビットのバイナリ数に変換することができる。すなわち、例えば、文字形状コード“A”はバイナリ数“000”に変換され、“U”は“001”に変換され、文字形状コード“g”は“010”に変換される。この限定セットの他の文字形状コードは、同様に、ユニークな3ビットのバイナリ数に変換される。トークン化テキスト原稿が、ステップS1010において、文字形状コードの限定セットからバイナリ数の列に変換されると、このビット列は、ステップS1020において、バイトの境界に渡って覆う3ビットのバイナリ・コードを必要なものとして有する、8ビットのバイトにグループ化される。すなわち、 $(2 + 2/3)$ 個の形状コードが各8ビットのバイトに圧縮される。続いて、ステップS1030において、圧縮された「コンパクト化されたトークン化」コードは、[コンパクト化コードのバイト数][コンパクト化コードのバイト]の形で記憶される。続いて、処理はステップS1040で終了する。

【0040】図8は、検索文字列(検索ストリング)を使用してコンパクト化されたトークン化テキスト・ファイルを検索する方法の第1の好ましい実施例を示している。ステップS1100からスタートし、ステップS1110において、検索文字列が入力される。ステップS1120において、検索文字列は、ステップS310～S530およびステップS1000～S1040の方法を使用して、まず、トークン化され、続いて、コンパクト化される。次に、ステップS1130において、コンパクト化されたトークン化検索文字列は、コンパクト化されたトークン化データベース内において、第1のコンパクト化されたトークン化原稿と比較され、検索文字列と検索している原稿の一部との間で正確に一致するものが突き止められる。ステップS1140において、システムは、正確に一致するものがコンパクト化されたトークン化検索文字列とコンパクト化された

トークン化テキスト原稿との間に存在するかどうかを判定する。ステップS1150において、正確に一致するものがステップS1140で突き止められると、制御はステップS1160に進む。ステップS1160で、一致ファイルと正確に一致するものを有するコンパクト化されたトークン化テキスト原稿の名前が出力される。続いて、制御はステップS1170に進む。一方、ステップS1150において、正確に一致するものがステップS1140で発見されなかったならば、制御はステップS1170に直接進む。

【0041】ステップS1170において、システムはデータベースの全てのエントリが分析されたかどうかを判定する。そうでなければ、制御はステップS1180に進む。ステップS1180において、次の原稿が選択される。続いて、制御はステップS1130に戻る。データベース内の全ての原稿が検索されると、制御はステップS1190に進み、一致するファイルを出力し、テキスト文字列と一致したデータベース・ファイルをリストする。続いて、制御はステップS1200に進み、処理は終了する。

【0042】このようにして、第1の好ましい実施例（正確に一致するものを発見する）においては、テスト文字列は、直接入力されるか、またはテキスト画像から取り出される。テスト文字列は、トークン化され、かつ、コンパクト化される。続いて、テスト文字列は、コントローラ114のコンパレータ36を用いて、コンパクト化されたトークン化テキスト画像またはコンパクト化されたトークン化電子テスト原稿の各エントリとそれぞれ比較され、正確に一致するものが突き止められる。正確に一致するものを有する各テキスト原稿は、テキスト文字列を有するデータベースの各原稿を識別するために出力される。

【0043】処理において、コンパクト化されたトークン化テスト文字列は、コンパクト化されたトークン化原稿または画像と、左から右へビットごとに、各原稿のビット列の開始部からビット列の終了部に向けて比較される。

【0044】第1の好ましい実施例の第1の変形例においては、ステップS1110～S1140でテスト文字列の入力および比較を行うのではなく、検索原稿全体が使用される。したがって、検索原稿とデータベースの複数の原稿のそれぞれとの比較の準備のために、原稿全体がトークン化され、かつ、コンパクト化される。

【0045】第1の好ましい実施例の第2の変形例においては、分類手段32は、最初に結合コンポーネントをステップS430において“x”の文字形状コード・シンボルに変換した後に、結合コンポーネントの右側の中央部分がベースライン位置とxラインの位置との間で開いた部分を有するかどうかを判定する。そのような開きがある場合には、文字分類手段34は文字形状コード“x”を文字形状コード“e”に変換する。文字形状コード“e”は、文字“e”および“c”を代表するものである。この第1の実施例の第2の変形例において、トークン化画像またはトークン化原稿がステップS1010で縮小されると、限定されたコード・セットは文字形状コード[CR]ではなくむしろ文字形状コード“e”を含むものとなる。一方、この第1の実施例の第2の変形例では、一致処理は、比較ステップの実行のために、第1および第2の好ましい実施例のいずれも使用することができる。

【0046】第2の好ましい実施例において、電子検索文字列は、トークン化されたテキスト画像のデータベースを検索するために使用される。

【0047】第3の好ましい実施例において、コンパクト化手段34は、ステップS1000～S1040のように、トークン化原稿を縮小トークン化原稿に変換して、この原稿をコンパクト化する代わりに、縮小トークン化原稿を使用して、少なくとも一つの浮動小数点数のセット（集合）を生成する。表1に示すように、少なくとも5つの可能な浮動小数点数を、値1～5として識別し、生成することができる。

【0048】

【表1】

値	定 義
1	素数をマップされた連続した形状コードの全てのペアをXOR 演算し、その結果を加算したもの
2	「平均形状コード」(各文字セルの形状コードの素数値の合計を形状コードの個数で割ったもの)
3	マップされた素数によって乗算されたときの各形状コードの発生の合計
4	連続したコードのペアをXOR 演算する代わりに、全ての他のコードがXOR 演算される点を除いて、値1と同じもの
5	ワード境界がクロスしない点を除いて値1と同じもの。すなわち、所与 “of the” $\text{map}[o] \wedge \text{map}[0], \text{map}[f] \wedge \text{map}[o], \text{map}[t] \wedge \neg \text{map}[o], \text{map}[h] \wedge \text{map}[t],$ $\text{map}[e] \wedge \text{map}[h]$ のXOR

【0049】表1に示すように、第1の値は、各文字形状コード・シンボルを数字のセットの一つに変換することにより生成される。数字のセットを構成する数字は、すべて相対的に互いに素である。続いて、連続した形状コードの各ペアにXOR 演算が施される。各XOR 演算の結果は加算され、第1の値が生成される。

【0050】第2の値は、トークン化ファイルの文字を表す数字を加算して、その合計値をテキスト文字列または原稿のトークン化文字の全体の個数によって割ることにより生成される。第3の値は、各文字形状コードがテキスト文字列または原稿に現れる回数をそれを表す数字に

乗じ、続いて各文字コードの乗算結果を加算することにより形成される。

【0051】第4の値は、各隣接したペアではなく、全ての他のコードがXOR 演算される点を除いて、第1の値と同様にして生成される。最後に、第5の値は第1の値と同様にして生成される。しかしながら、第5の値において、XOR 演算は、ワードの境界に渡って拡張されない(各ワード・トークンのリーディング・スペースを含むが)。

【0052】データベースの各テキスト画像または原稿に対するこれらの値の少なくとも一つと、検索文字列または原稿に対する対応する値とを生成することによって、テキスト文字列または原稿とデータベースを形成する画像またはテキスト原稿との間の正確でない(厳密でない)一致を突き止めることができる。すなわち、図9に示すように、ステップS2000 でスタートした後、テスト原稿の少なくとも一つの浮動小数点数のセットがステップS2010 で生成される。続いて、ステップS2020 において、テスト原稿の浮動小数点数のセットのそれぞれ一つの値は、原稿データベースのテキスト画像またはテキ

スト原稿の現在のものの浮動小数点数のセットのそれぞれ一つの対応する浮動小数点数と比較される。続いて、ステップS2030 において、距離メジャーが、原稿データベースのエントリの現在のもの用に生成される。この距離メジャーは、テスト原稿の浮動小数点の値のセットと原稿データベースの現在のエントリの浮動小数点の値のセットとの間の全体の相違および類似を表す。

【0053】続いて、この距離の値はステップS2040 で分析され、テスト原稿と現在のエントリとの間の少なくとも厳密でない一致点を示しているかどうかが判定される。ステップS2050 で、少なくとも厳密でない一致点が発見されると、ステップS2060 において、一致する原稿の名前が一致ファイルへ出力される。一致するものが発見されないならば、制御はステップS2050 から直接ステップS2070 に進む。続いて、ステップS2070 において、データベースはチェックされ、全てのエントリがテスト原稿と比較されたかどうか判断される。そうでなければ、制御はステップS2080 に進む。ステップS2080 において、次のエントリが選択され、制御はステップS2030 に戻る。そうであれば、少なくとも厳密でない一致するエントリのリストがステップS2090 で出力される。ステップS2100 でシステムは停止する。検索文字列またはデータベースとテキスト画像またはテキスト原稿との間の少なくとも厳密でない一致点として識別されたものが正当であることを確認するために、しきい値が用いられる。これに加えて、オペレータがテキスト画像またはテキスト原稿と検索原稿とをどの程度の近似度で一致させたいかに依存して、このしきい値は調整可能である。

【0054】上述した第1の好ましい実施例を用いて、英語、フランス語およびドイツ語の原稿を有する多重言語原稿データベースが、これらの言語の話し手に興味の

ある167個の異なるテスト文字列を用いて検索された。UNIXの“AGREP”コマンド（近似した文字列の一致検索を行う）を用いて、文字列のコンパクト化され、トークン化された形態のものがデータベースの原稿のコンパクト化され、トークン化された形態のものと比較された。167個の検索文字列の68%（114個の検索文字列）が、フォールス・ポジティブが返されず、すなわちフォールス・ネガティブが返された。他の32%（54個の検索文字列）は、エラー・レートが増加した。しかしながら、これらの54個の検索文字列は、より短く、あまり特異性のない文字列になる傾向があった。

【0055】第2の好ましい実施例において、46個の英語、フランス語およびドイツ語の原稿がトークン化され、かつ、コンパクト化されて、データベースを形成している。これらの原稿のうちの36個の走査されたものが、走査された原稿をトークン化するが、データベースの原稿をコンパクト化しないことにより、オリジナル原稿と比較された。続いて、走査された原稿の縮小トークン化形態が、オリジナル原稿の縮小トークン化形態と比較された。36個のオリジナル原稿を生成するとき、特別な注意は払われなかった。したがって、これらの原稿は1.8°上向きにスキューしていた。このテストにおいて、36個の全原稿はそのオリジナルと正確に一致した。さらなるテストにおいて、写真複写の複数の生成物が走査され、データベース原稿と比較された。この第2の好ましい実施例は、3番目および4番目の写真複写生成物にまで耐えた。さらに、この第2の好ましい実施例は、場合によっては、7番目の写真複写生成物でさえも正確に一致させることができた。

【0056】これらの上記テストは300dpiの走査解像度で行われた。さらなるテストは、200dpiのテキスト画像を走査することにより行われた。このより低い解像度のテストにおいて、36個の原稿のうちの34個の原稿が、データベースのオリジナルの原稿と正確に一致した。2つの原稿がデータベースの原稿のいずれにも一致しなかったが、フォールス・ポジティブは生成されなかった。

【0057】最後のテストは浮動小数点のフルの5次元

セットを使用し、これらの36個の走査された原稿がデータベースの原稿と比較された。第4の好ましい実施例が最も近似して一致するものを発見するために試験され、一致度95%のしきい値が使用された。この第4の好ましい実施例を用いて、36個の走査されたテキスト画像のうちの34個の対応するオリジナル原稿が最も近似して一致するものとして識別された。他の2つのものでは、対応するオリジナル原稿は2番目に近似して一致するものであった。

【図面の簡単な説明】

【図1】文字形状コード認識システムのブロック図である。

【図2】オリジナルの原稿を示す。

【図3】文字形状コードのシンボルおよび実際の対応するスクリプト文字を示す。

【図4】本発明による文字形状コード分類方法の好ましい実施例のフローチャートを示す。

【図5】文字形状コード分類の決定木を示す。

【図6】図3の文字形状コード分類計画によって変換された図2のテキスト部分を示す。

【図7】コンパクト化方法の第1の好ましい実施例のフローチャートを示す。

【図8】テキスト・ファイルの検索およびマッチング方法の第1の好ましい実施例のフローチャートを示す。

【図9】テキスト・ファイルの検索およびマッチング方法の第4の好ましい実施例のフローチャートを示す。

【符号の説明】

110 スキャナ

112 メモリ

114 コントローラ

24 入出力装置

26 画像メモリ

28 結合コンポーネント生成器

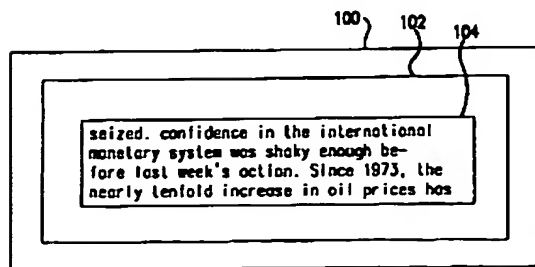
30 空間的特徴決定手段

32 文字タイプ分類手段

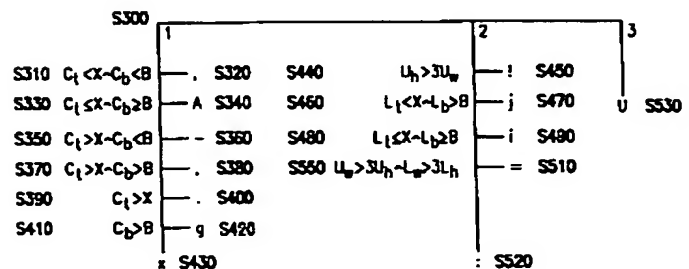
34 コンパクト化手段

36 比較器

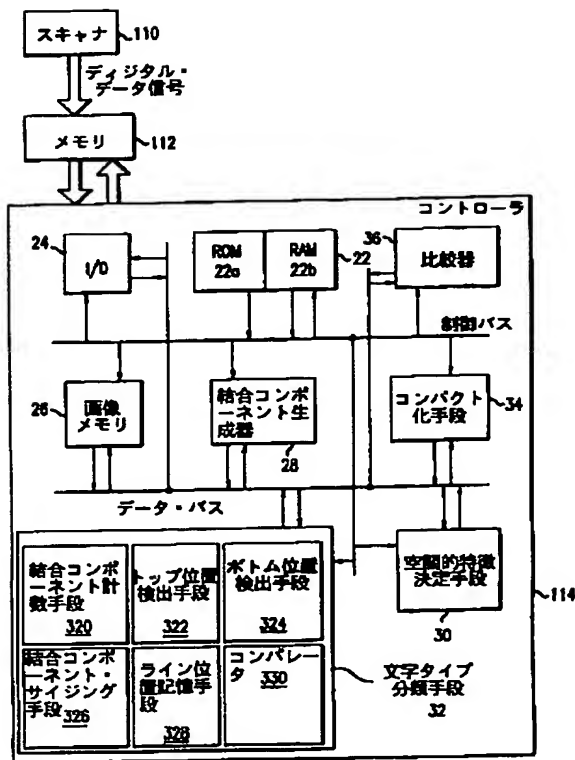
【図2】



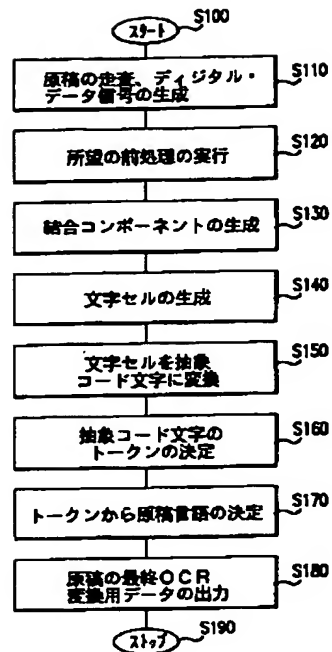
【図5】



【図 1】



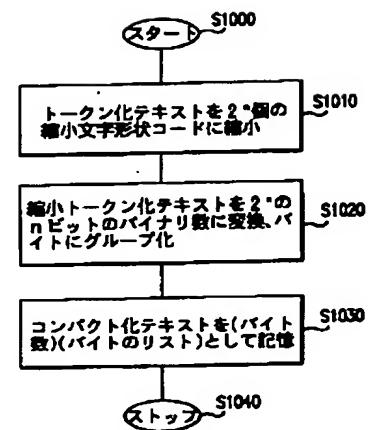
【図 4】



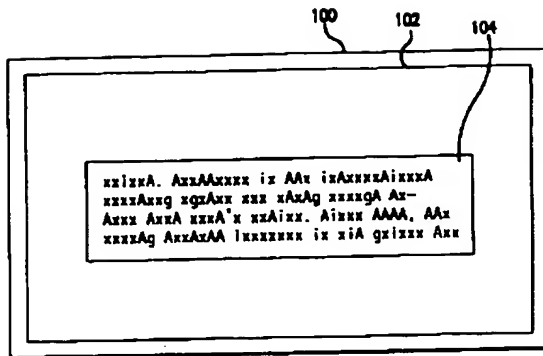
【図 3】

文字形状コード	メンバー
A	ABCDEFGHIJKLMNOPQRSTUVWXYZabcdefghijklmnopqrstuvwxyz01234567890a0101H
x	abcdefghijklmnopqrstuvwxyz
i	abcdefghijklmnopqrstuvwxyz
g	8P977
j	j
.	.
-	-
.	.
.	.
i	i
e	e
U	abcdefghijklmnopqrstuvwxyz
!	!

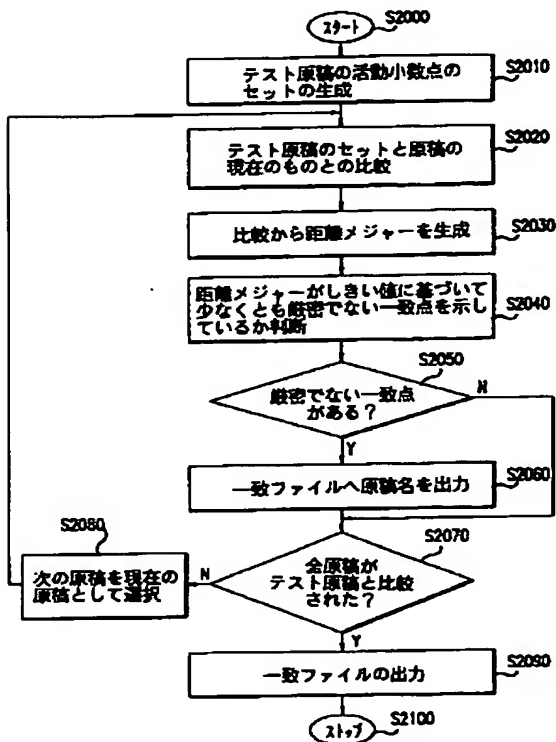
【図 7】



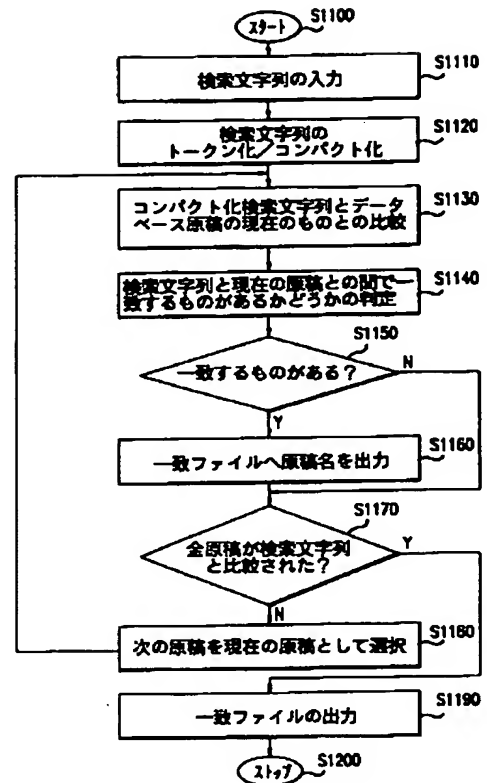
【図 6】



【図 9】



【図 8】



フロントページの続き

(72) 発明者 アンтониオ ビー ディアス
 アメリカ合衆国 マサチューセッツ州
 02138ケンブリッジ ハーヴァード ユニ
 ヴァーシティー クインシー メール セ
 ンター 200